

Posterior Mean Super-resolution with a Causal Gaussian Markov Random Field Prior

Takayuki Katsuki, Akira Torii, and Masato Inoue

Abstract—We propose a Bayesian image super-resolution (SR) method with a causal Gaussian Markov random field (MRF) prior. SR is a technique to estimate a spatially high-resolution image from given multiple low-resolution images. An MRF model with the line process supplies a preferable prior for natural images with edges. We improve the existing image transformation model, the compound MRF model, and its hyperparameter prior model. We also derive the optimal estimator – not the joint maximum a posteriori (MAP) or marginalized maximum likelihood (ML), but the posterior mean (PM) – from the objective function of the L2-norm (mean square error) -based peak signal-to-noise ratio (PSNR). Point estimates such as MAP and ML are generally not stable in ill-posed high-dimensional problems because of overfitting, while PM is a stable estimator because all the parameters in the model are evaluated as distributions. The estimator is numerically determined by using variational Bayes. Variational Bayes is a widely used method that approximately determines a complicated posterior distribution, but it is generally hard to use because it needs the conjugate prior. We solve this problem with simple Taylor approximations. Experimental results have shown that the proposed method is more accurate or comparable to existing methods.

Index Terms—super-resolution, Bayesian inference, Markov random field prior, line process, posterior mean, variational Bayes, Taylor approximation.

I. INTRODUCTION

Super-resolution (SR) is an information processing technique that makes it possible to infer a spatially high-resolution (HR) image of a scene from corresponding multiple low-resolution (LR) images that are affected by warping, blurring, and noise. SR can be applied to a variety of images; e.g., still images extracted from several sequential video frames. SR needs the registration of LR images in addition to the image restoration of the registered LR images. Since the earliest work by Tsai and Huang [1], SR has been achieved using various methods [2]–[10] and good overviews of these methods are given in [11]–[16]. Generally, SR is an ill-posed inverse problem because inverting the blur process without amplifying the effect of the noise is difficult [13]. In other words, the degrees of freedom of the HR image and pixel-wise observation noise are always higher than the dimensionality of the observed LR images, so complete determination of an HR image is impossible. Therefore, the HR image is frequently inferred as the most preferable image within the framework of the probabilistic information processing, and we handle SR using this framework in this paper. The probabilistic

information processing has three key features: 1) model, 2) objective function, and 3) optimization method. In the SR problem, the model includes the observation model and the prior model. The observation model consists of warping, blurring, downsampling, and noise models. The prior model, necessary for the Bayesian framework, mainly consists of an HR image prior, and sometimes includes both the hyperparameter prior for the HR image prior and the registration prior. The objective function evaluates how good or bad an estimator is. The estimator usually represents the inferred HR image, and sometimes includes auxiliary parameters; e.g., the registration parameters and edge information. The optimization method numerically maximizes/minimizes the objective function and determines the estimator. An optimization method is not necessary for simple problems in which an analytical exact solution can be obtained. In the probabilistic information processing, SR can be categorized according to these three key features.

To deal with warping, blurring, and downsampling, a linear transformation model is frequently used [3], [6], [8], [10]. Warping is usually limited with planar rotation and parallel translation. Blurring is defined by using a point spread function (PSF); a square or Gaussian type PSF is common. Downsampling denotes sampling from an HR image to construct an LR image. Downsampling sometimes includes anti-aliasing. Since these three transformations are linear, they can be combined into a single transformation matrix. As for the noise model, pixel-independent additive white Gaussian noise (AWGN) is usually employed.

The Bayesian framework, especially the HR image prior, is quite useful for SR. The HR image prior provides appropriate smoothness between neighboring pixel luminances. A common type of HR image prior imposes an L2-norm penalty on differences between horizontally and vertically adjacent pixel luminances (the first derivative). The L1-norm of the first derivative is sometimes used, and it has the advantage of robust inference against outliers. The total variation (TV) prior [10] employs the L1-norm of the gradient vector. The Huber prior [5] is a mixture prior of L1- and L2-norms. The SAR model [2], [9], [17] employs the response of a two-dimensional Laplacian filter (the second derivative). The Gaussian process prior [3] has neighboring pixels spread according to a Gaussian distribution. Besides the degree of smoothness between neighboring pixels, information regarding the discontinuity, or equivalently, the edges or line process, is also useful for inference. A common type of prior implementing edges is the compound Markov random field (MRF) prior that was introduced by Geman & Geman [18] and is widely used [4], [6], [8]. With respect to the compound MRF

Takayuki Katsuki, Akira Torii, and Masato Inoue are with the Department of Electrical Engineering and Bioscience, Graduate School of Advanced Science and Engineering, Waseda University, 3-4-1, Okubo, Shinjuku, Tokyo 1698555, Japan. E-mail: (see http://www.eb.waseda.ac.jp/m_inoue/).

[19], [20] prior, the normalizing constant, or equivalently, the partition function, is usually difficult to calculate because it has an exponential calculation cost with respect to the dimensionality of the line process. Recently, Kanemura *et al.* [6], [8] confusingly introduced a “causal” type of Gaussian MRF prior whose calculation cost is polynomial. We try to improve this prior in this paper.

The SR estimator should be derived from an objective function. As the objective function, a posterior distribution has been widely employed. Since the posterior distribution usually includes both the HR image and registration parameters, the joint maximum a posteriori (MAP) solution [2] is a suitable estimator for this objective function. Other than the joint MAP, the use of the marginalized maximum likelihood (ML) [3], [6] or marginalized MAP [5] has been proposed. Tipping *et al.* [3] and Kanemura *et al.* [6], [8] determine the registration parameters by using ML inference, where the HR image is marginalized out, and determine the HR image by using MAP inference. Pickup *et al.* [5] determines the HR image by using MAP inference, wherein the registration uncertainties are marginalized out, and assumes that the registration parameters are pre-registered by using standard registration techniques. Marginalized ML is also called type-II ML, evidence approximation, or empirical Bayes. Marginalized ML has no registration prior, unlike marginalized MAP. Pickup *et al.* [5] reported that marginalized MAP is superior to both joint MAP and marginalized ML. We evaluate the accuracy of SR methods in terms of the L2-norm (mean square error) - based peak signal-to-noise ratio (PSNR). Therefore, we think it is natural to employ PSNR as the objective function. For this objective function, posterior mean (PM) is a suitable estimator. The variational Bayes [21] approach [10] seems to approximately determine the PM of the HR image, although the authors assume some registration parameters are known and use point-estimate model parameters obtained by ML inference. To determine the exact PM of the HR image, all parameters other than the HR image should be marginalized out over the joint posterior distribution.

The type of optimization method to use is not as substantial a problem as the choice of model and objective function, but it is still important. Since almost all good estimators cannot be exactly determined because of difficult analytical integration or an exponential calculation cost, some approximation methods need to be introduced. Also, parameter tuning is necessary in many numerical optimization methods; e.g., of the initial value and the step-width settings in gradient methods. Specifically, in early work done on image restoration, an annealing method was used for the joint MAP solution [18], [22]. For marginalized ML and marginalized MAP solutions, the scaled conjugate gradients algorithm was used [3], [5]. In recent work, the variational expectation-maximization (EM) algorithm has been applied, which includes the gradient method in the M step [6], [8]. The variational Bayes approach has also been applied [10]. This method includes nested optimization of the majorization-minimization approach. This majorization-minimization approach seems to affect both the HR image prior and the estimator. Specifically, it modifies the TV prior to include a discontinuity parameter (called local spatial activity).

In addition, this parameter is point-estimated when the HR image is inferred.

In this paper, we propose a new SR method that employs a “causal” Gaussian MRF prior and utilizes variational Bayes to calculate the optimal estimator, PM, with respect to the objective function of the L2-norm-based PSNR. This is a straightforward approach, but it was not proposed earlier possibly because an important limitation of variational Bayes is that a conjugate prior is needed. We solve this problem through simple Taylor approximations. In Section II, we define models, where we introduce a novel unified warping, blurring and downsampling model, an improved HR image prior, an improved hyperparameter prior, and a registration prior. In Section III, we employ PSNR as the objective function and derive the optimal estimator, PM, from this objective function. In Section IV, we determine the PM by using variational Bayes and Taylor approximations. In Section V, we evaluate the proposed method by comparing it with existing methods. We discuss the proposed method in Section VI and conclude in Section VII.

II. MODEL

A. Definitions

First, we define the gamma, Bernoulli, and Gaussian distributions used in this paper:

$$\text{Gamma}(x; a, b) \equiv \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx} \quad (x > 0),$$

$$\text{Bernoulli}(x; \mu) \equiv \mu^x (1 - \mu)^{1-x} \quad (x \in \{0, 1\}),$$

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (\mathbf{x} \in \mathbb{R}^d),$$

Here, Γ is the gamma function, $|\bullet|$ denotes the determinant of a given matrix, superscript \top denotes the transpose, \mathbb{R} is the real number field, and d is the dimension of \mathbf{x} . The logistic function and Kullback-Leibler (KL) divergence from distributions $p(\mathbf{x})$ to $q(\mathbf{x})$ are respectively defined as

$$\text{logistic}(x) \equiv \frac{1}{1 + e^{-x}},$$

$$D_{\text{KL}}(p(\mathbf{x}) \| q(\mathbf{x})) \equiv \left\langle \ln \frac{p(\mathbf{x})}{q(\mathbf{x})} \right\rangle_{p(\mathbf{x})},$$

where the angle brackets $\langle \bullet \rangle_{\circ}$ denote the expectation of \bullet with respect to a distribution \circ . Additionally, tr denotes the trace of a given matrix. diag denotes a diagonal matrix. \mathbf{I} is an identity matrix of appropriate size. $\mathbf{0}$ is a zero vector or a zero matrix of appropriate size. All the vectors in this paper are column vectors. The $\|\bullet\|_2$ denotes the L2-norm of a given vector. At this point, these variables have absolutely nothing to do with the variables that appear later.

B. Observation Model

Our task is to estimate an HR grayscale image, $\mathbf{x} \in \mathbb{R}^{N_x}$, from the observed multiple LR grayscale images, $\mathbf{Y} \equiv \{\mathbf{y}_l\}_{l=1}^L$, $\mathbf{y}_l \in \mathbb{R}^{N_y}$. Images \mathbf{y}_l and \mathbf{x} are regarded as lexicographically stacked vectors. The number of pixels for each LR image, N_y , is assumed to be less than that of the HR image, N_x ; i.e., $N_y < N_x$. We do this estimation using

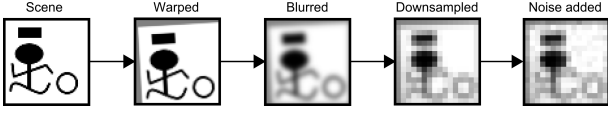


Fig. 1. An illustration of the image observation process

an SR technique whose resolution enhancement factor is $\alpha \equiv \sqrt{N_x/N_y} (> 1)$. Although we define the range of a pixel luminance value as infinite, we use -1 for black, $+1$ for white, and values between -1 and $+1$ for gradual gray.

The image observation process is modeled as shown in Fig. 1; the HR image \mathbf{x} is geometrically warped, blurred, downsampled, and corrupted by noise ϵ_l to form the observed LR image \mathbf{y}_l :

$$\mathbf{y}_l \equiv \mathbf{W}(\phi_l)\mathbf{x} + \epsilon_l, \quad (1)$$

or, more strictly,

$$p(\mathbf{Y}|\mathbf{x}, \beta, \Phi) \equiv \prod_{l=1}^L \mathcal{N}(\mathbf{y}_l; \mathbf{W}(\phi_l)\mathbf{x}, \beta^{-1}\mathbf{I}). \quad (2)$$

The $\epsilon_l \in \mathbb{R}^{N_y}$ is AWGN with precision (inverse variance) $\beta (> 0)$. Here, $\mathbf{W}(\phi_l)$ is the $N_y \times N_x$ transformation matrix that is simultaneously used for warping, blurring, and downsampling. It is defined as

$$\mathbf{W}(\phi_l)_{j,i} \equiv \frac{\mathcal{N}(\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i); 0, \gamma_l^{-1}\mathbf{I})}{\sum_{i' \in \mathcal{I}} \mathcal{N}(\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_{i'}); 0, \gamma_l^{-1}\mathbf{I})}, \quad (3)$$

$$\vec{\chi}(\theta, \vec{o}, \vec{\zeta}, \vec{\xi}) \equiv \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} (\alpha \vec{\zeta} - \vec{o}) - \vec{\xi}, \quad (4)$$

where \mathcal{I} represents the extent of the summation (explained in the next paragraph), and the vectors $\vec{\xi}_i$ and $\vec{\zeta}_j$ respectively denote the two-dimensional positions of the i -th pixel of the original HR image and the j -th pixel of the observed LR image. We define the center of each image as the origin and the size of each pixel is 1 by 1. For example, regarding an HR image with 40×40 pixels, each ξ represents $[-19.5, -19.5]^\top, [-18.5, -19.5]^\top, \dots, [19.5, 19.5]^\top$. θ_l and \vec{o}_l represent the warping parameters of the l -th LR image: the rotational motion parameter and translational motion parameter. The Gaussian distribution in (3) represents a Gaussian PSF that defines the blur, and $\gamma_l (> 0)$ represents its precision parameter. In this paper, we assume γ_l also differs for each observed image. These transformation parameters are packed into ϕ_l , which is defined as

$$\Phi \equiv \{\phi_l\}_{l=1}^L, \quad \phi_l \equiv [\phi_l, k]_{k=1}^4 \equiv [\theta_l, [\vec{o}_l]_h, [\vec{o}_l]_v, \gamma_l]^\top, \quad (5)$$

where subscripts h and v , respectively, denote horizontal and vertical positions on the image.

In previous works [3], [6], [8], the extent of \mathcal{I} was defined as the extent of the HR image. According to this definition, however, the shape of the PSF is no longer Gaussian. For example, at the corner of the HR image, the shape is not omnidirectional but limited in a way such as that of a quadrant. In this paper, the extent of \mathcal{I} is defined as infinite, and the luminance values outside the HR image are defined as 0

(middle gray). This normalization term faithfully represents the Gaussian PSF. We also found that this normalization term is exactly given by using the elliptic theta function ϑ_3 , and we can rewrite $\mathbf{W}(\phi_l)$ as

$$\begin{aligned} \mathbf{W}(\phi_l)_{j,i} &= \frac{\mathcal{N}(\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i); 0, \gamma_l^{-1}\mathbf{I})}{\vartheta_3\left(\left[\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i)\right]_h, e^{-\frac{2\pi^2}{\gamma_l}}\right)\vartheta_3\left(\left[\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i)\right]_v, e^{-\frac{2\pi^2}{\gamma_l}}\right)}, \end{aligned} \quad (6)$$

$$\vartheta_3(u, q) \equiv 1 + 2 \sum_{n=1}^{\infty} q^{n^2} \cos 2n\pi u. \quad (7)$$

The elliptic theta function includes an infinite series, but it is easily determined numerically because the convergence is quite fast. In (6), the normalization term (the denominator of the right-hand side) seems to depend on i because $\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i)$ includes $\vec{\xi}_i$, but this is not true. Because the elliptic theta function is a periodic function with respect to the argument u with period 1, and $\vec{\chi}(\theta_l, \vec{o}_l, \vec{\zeta}_j, \vec{\xi}_i)$ can only take discrete values with step size 1 for the horizontal and vertical directions, the normalization term has the same value with respect to i .

C. HR Image Prior

Here, we introduce a “causal” Gaussian MRF prior for the HR image and additional latent variables. These latent variables are called the line process that controls the local correlation among pixel luminances. The introduction of the latent variables enables explicit expression of the possible discontinuity in the HR image. The line process, $\boldsymbol{\eta}$, consists of binary variables $\eta_{i,j} \in \{0, 1\}$ for all adjacent pixel pairs i and j . Its size equals $N_\eta \equiv 2N_x - [\text{number of HR image's horizontal pixels}] - [\text{number of HR image's vertical pixels}]$. We define the prior as

$$p(\mathbf{x}, \boldsymbol{\eta} | \lambda, \rho, \kappa) \equiv p(\mathbf{x} | \boldsymbol{\eta}, \rho, \kappa) p(\boldsymbol{\eta} | \lambda) \quad (8)$$

$$\begin{aligned} &= \exp \left[-\lambda \sum_{i \sim j} (1 - \eta_{i,j}) - \frac{\rho}{2} \sum_{i \sim j} \eta_{i,j} (x_i - x_j)^2 - \frac{\kappa}{2} \|\mathbf{x}\|_2^2 \right. \\ &\quad \left. + \frac{1}{2} \ln \left| \frac{\mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)}{2\pi} \right| + N_\eta \ln \text{logistic}(\lambda) \right], \end{aligned} \quad (9)$$

where

$$p(\boldsymbol{\eta} | \lambda) \equiv \prod_{i \sim j} \text{Bernoulli}(\eta_{i,j}; \text{logistic}(\lambda)), \quad (10)$$

$$p(\mathbf{x} | \boldsymbol{\eta}, \rho, \kappa) \equiv \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)^{-1}), \quad (11)$$

$$\mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)_{i,j} \equiv \begin{cases} \rho \sum_{k \sim i} \eta_{i,k} + \kappa, & i = j, \\ -\rho \eta_{i,j}, & i \sim j, \\ 0, & \text{otherwise.} \end{cases} \quad (12)$$

Here, the summation $\sum_{i \sim j}$ is taken over all pairs of adjacent pixels. The notation $i \sim j$ means that the i -th and j -th pixels are adjacent in the upward, downward, leftward, and rightward directions. The line process $\boldsymbol{\eta}$ switches the local characteristics

of the prior. It indicates whether two adjacent pixels take similar values or independent values. When $\eta_{i,j} = 1$, the i -th and the j -th pixels are strongly smoothed according to the quadratic penalty, whereas there is no smoothing when $\eta_{i,j} = 0$. The hyperparameter $\lambda (> 0)$ is an edge penalty parameter that prevents $\eta_{i,j}$ from excessively taking edges. Note that λ is restricted to positive values because a negative λ leads to a reward rather than a penalty for taking edges. $\rho (> 0)$ is a smoothness parameter that prevents the differences in adjacent pixel luminances from becoming large, and $\kappa (> 0)$ is a contrast parameter that prevents \mathbf{x} from taking an improperly large absolute value. On the other hand, in previous works [6], [8], κ is assumed to be 0, which results in an improper normalizing constant (see Discussion). $\mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)$ is the $N_{\mathbf{x}} \times N_{\mathbf{x}}$ precision matrix of \mathbf{x} .

We have defined the introduced causal Gaussian MRF prior in the joint distribution form of \mathbf{x} and $\boldsymbol{\eta}$, i.e., $p(\boldsymbol{\eta})p(\mathbf{x}|\boldsymbol{\eta})$. We call such a model “causal” because $\boldsymbol{\eta}$ seems to cause \mathbf{x} . The MRF model is defined as having the property

$$p(x_i|\mathbf{x}\setminus x_i, \boldsymbol{\eta}) = p(x_i|\mathbf{x}_{\mathcal{L}(i)}, \boldsymbol{\eta}_{i,\mathcal{L}(i)}) \quad (13)$$

in this case; i.e., the conditional distribution of a random variable, x_i , given all other variables, $\mathbf{x}\setminus x_i$ and $\boldsymbol{\eta}$, equals the conditional distribution of the random variable, x_i , given its “neighboring” variables, $\mathbf{x}_{\mathcal{L}(i)}$ and $\boldsymbol{\eta}_{i,\mathcal{L}(i)}$. If this conditional distribution is a Gaussian distribution, such an MRF is called a Gaussian MRF.

The “compound” MRF prior is usually defined in the form of the Gibbs distribution [18],

$$\tilde{p}(\mathbf{x}, \boldsymbol{\eta}) \equiv \frac{\exp(-\tilde{H}(\mathbf{x}, \boldsymbol{\eta}))}{\sum_{\boldsymbol{\eta}} \int \exp(-\tilde{H}(\mathbf{x}, \boldsymbol{\eta})) d\mathbf{x}}, \quad (14)$$

which is based on some microstate energy function, or equivalently, a Hamiltonian, such as

$$\begin{aligned} \tilde{H}(\mathbf{x}, \boldsymbol{\eta}) \\ \equiv \lambda \sum_{i \sim j} (1 - \eta_{i,j}) + \frac{\rho}{2} \sum_{i \sim j} \eta_{i,j} (x_i - x_j)^2 + \frac{\kappa}{2} \|\mathbf{x}\|_2^2. \end{aligned} \quad (15)$$

In addition to the property of (13), a compound MRF also has the property of

$$\tilde{p}(\eta_{i,j}|\mathbf{x}, \boldsymbol{\eta}\setminus \eta_{i,j}) = \tilde{p}(\eta_{i,j}|x_i, x_j), \quad (16)$$

whereas the introduced “causal” Gaussian MRF prior does not. Therefore, we do not call the introduced prior a “compound” MRF prior, even though (8) and (14) have similar forms. Furthermore, the introduced “causal” Gaussian MRF prior is a generative model, whereas the “compound” MRF is not. A generative model has the advantage of reducing the calculation cost (see Discussion).

D. Hyperparameter Prior

Generally, prior distributions should be non-informative unless we have explicit reasons because an informative prior leads to heuristics. Actually, we define the prior distributions

for the hyperparameters of the HR image prior to be as non-informative as possible:

$$\begin{aligned} p(\lambda, \rho, \kappa, \beta) \equiv & \text{Gamma}(\lambda; a_{\lambda}^{(0)}, b_{\lambda}^{(0)}) \text{Gamma}(\rho; a_{\rho}^{(0)}, b_{\rho}^{(0)}) \\ & \times \text{Gamma}(\kappa; a_{\kappa}^{(0)}, b_{\kappa}^{(0)}) \text{Gamma}(\beta; a_{\beta}^{(0)}, b_{\beta}^{(0)}), \end{aligned} \quad (17)$$

$$\begin{aligned} a_{\lambda}^{(0)} \equiv 10^{-2}, b_{\lambda}^{(0)} \equiv 10^{-2}, a_{\rho}^{(0)} \equiv 10^{-2}, b_{\rho}^{(0)} \equiv 10^{-2}, \\ a_{\kappa}^{(0)} \equiv 10^{-2}, b_{\kappa}^{(0)} \equiv 10^{-2}, a_{\beta}^{(0)} \equiv 10^{-2}, b_{\beta}^{(0)} \equiv 10^{-2}. \end{aligned} \quad (18)$$

For a gamma distribution, the number of effective prior observations in the Bayesian framework is equal to two times parameter a . As shown in the Appendix, the number of observations for the hyperparameter λ is $N_{\boldsymbol{\eta}}$ in this SR. Also, that for ρ and κ is $N_{\mathbf{x}}$, and that for β is $LN_{\mathbf{y}}$. Therefore, the above settings – e.g., $2a_{\lambda}^{(0)} \ll N_{\boldsymbol{\eta}}$ – are considered sufficiently non-informative. Superscript (0) is added because we use these parameters as the initial values of variational Bayes later.

E. Registration Prior

For the registration parameters including the blurring parameter, we also define the corresponding prior as

$$\begin{aligned} p(\boldsymbol{\Phi}) \equiv & \prod_{l=1}^L \mathcal{N}(\phi_l; \boldsymbol{\mu}_{\phi_l}^{(0)}, \boldsymbol{\Sigma}_{\phi_l}^{(0)}), \\ \boldsymbol{\mu}_{\phi_l}^{(0)} \equiv & [0, 0, 0, 12/\alpha^2], \quad \boldsymbol{\Sigma}_{\phi_l}^{(0)} \equiv \text{diag}[10^{-3}, 10^0, 10^0, 10^{-3}]. \end{aligned} \quad (19)$$

For the rotational motion parameter θ_l , the prior assumes 0 ± 1.81 degree ($\frac{180}{\pi} \sqrt{10^{-3}} \simeq 1.81$). This assumption is considered suitable for this SR task. Similarly, an assumption of 0 ± 1 pixels for translational motion parameters $[\vec{o}_l]_h$ and $[\vec{o}_l]_v$ is considered suitable. For blurring parameter γ_l , $\mu_{\gamma_l}^{(0)}$ is taken to be the value equivalent to the anti-aliasing of the scale factor α .

III. OBJECTIVE FUNCTION AND ESTIMATOR

A. Peak Signal-to-Noise Ratio (PSNR)

First, we confirm that the joint distribution of all random variables can now be explicitly given as

$$p(\mathbf{Y}, \mathbf{z}) = p(\mathbf{Y}|\mathbf{x}, \beta, \boldsymbol{\Phi})p(\mathbf{x}, \boldsymbol{\eta}|\lambda, \rho, \kappa)p(\lambda, \rho, \kappa, \beta)p(\boldsymbol{\Phi}), \quad (21)$$

$$\mathbf{z} \equiv [\mathbf{x}, \boldsymbol{\eta}, [\lambda, \rho, \kappa, \beta], \boldsymbol{\Phi}], \quad (22)$$

Once the joint distribution is obtained, we can derive all the marginal and conditional distributions; e.g., the posterior distribution $p(\mathbf{z}|\mathbf{Y})$ and joint distribution of the HR and LR images $p(\mathbf{Y}, \mathbf{x})$.

One of the most commonly used evaluation functions of the inferred image is the L2-norm (mean square error) -based PSNR. It is defined as

$$\text{PSNR}(\hat{\mathbf{x}}; \mathbf{x}) \equiv 10 \log_{10} \frac{2^2}{\frac{1}{N_{\mathbf{x}}} \|\hat{\mathbf{x}} - \mathbf{x}\|_2^2}, \quad (23)$$

where $\hat{\mathbf{x}}$ is the estimator of the HR image and \mathbf{x} is the true HR image. Since only LR images, \mathbf{Y} , are available for the estimator, we sometimes explicitly express it as a function

form, $\hat{\mathbf{x}}(\mathbf{Y})$. Now, our objective function (functional) to be maximized regarding the estimator is defined as

$$F(\hat{\mathbf{x}}(\mathbf{Y})) \equiv 10 \log_{10} \frac{2^2}{\left\langle \frac{1}{N_{\mathbf{x}}} \|\hat{\mathbf{x}}(\mathbf{Y}) - \mathbf{x}\|_2^2 \right\rangle_{p(\mathbf{Y}, \mathbf{x})}}. \quad (24)$$

This is because we prefer good estimator performance on average over various HR images and the corresponding LR images. Here, we assume that the occurrence rate of HR and LR images exactly coincides with the model we just introduced.

B. Posterior Mean (PM)

Using the above objective function, we can explicitly derive the best estimator of the HR image as the PM,

$$\operatorname{argmax}_{\hat{\mathbf{x}}(\mathbf{Y})} F(\hat{\mathbf{x}}(\mathbf{Y})) = \langle \mathbf{x} \rangle_{p(\mathbf{x}|\mathbf{Y})}. \quad (25)$$

Here, we used the well-known fact that the PM coincides with the minimum mean square error estimator in Bayesian framework. Note that $p(\mathbf{x}|\mathbf{Y})$ needs marginalization of all parameters other than \mathbf{x} over $p(\mathbf{z}|\mathbf{Y})$. If the PM of the line process or other model parameters is necessary, it can also be determined in the same manner.

IV. OPTIMIZATION METHOD

A. Variational Bayes

Though we could derive the optimal estimator, we cannot obtain the analytical solutions of the posterior distribution $p(\mathbf{z}|\mathbf{Y})$ and marginalized posterior distribution $p(\mathbf{x}|\mathbf{Y})$. Consequently, we have to rely on approximations. Here, we employ variational Bayes.

Variational Bayes [21] provides a trial distribution $q(\mathbf{z})$ that approximates the true posterior. We impose a factorization assumption on the trial distribution,

$$q(\mathbf{z}) \equiv q(\mathbf{x})q(\boldsymbol{\eta})q(\lambda, \rho, \kappa, \beta)q(\boldsymbol{\Phi}). \quad (26)$$

Note that, at this moment, the distribution family of each factorized distribution is not limited. We identify the optimal trial distribution that minimizes the KL divergence between the trial and the true distributions as the best approximation of the true distribution:

$$\hat{q}(\mathbf{z}) \equiv \operatorname{argmin}_{q(\mathbf{z})} D_{\text{KL}}(q(\mathbf{z}) \| p(\mathbf{z}|\mathbf{Y})). \quad (27)$$

Actually, the trial distribution that minimizes the KL divergence, not from $q(\mathbf{z})$ to $p(\mathbf{z}|\mathbf{Y})$ but from $p(\mathbf{z}|\mathbf{Y})$ to $q(\mathbf{z})$ coincides with the product of the exact marginal distributions as

$$\operatorname{argmin}_{q(\mathbf{z})} D_{\text{KL}}(p(\mathbf{z}|\mathbf{Y}) \| q(\mathbf{z})) = \prod_i p(z_i|\mathbf{Y}), \quad (28)$$

but this minimization is difficult to calculate.

Under the factorization assumption of the trial distribution and the extremal condition of the KL divergence, each optimal trial distribution should satisfy the self-consistent equations,

$$\hat{q}(z_i) \propto \exp(\ln p(z_i|\mathbf{Y})) \prod_{j \neq i} \hat{q}(z_j). \quad (29)$$

In the common style of variational Bayes [10], [23], this equation is solved by making repetitive updates,

$$q^{(0)}(z_i) \equiv p(z_i), \quad (30)$$

$$q^{(t+1)}(z_i) \propto \exp(\ln p(\mathbf{z}|\mathbf{Y})) \prod_{j \neq i} q^{(t)}(z_j). \quad (31)$$

Each factorized trial distribution is supposed to converge to the optimal distribution. Sometimes, some $q^{(t+1)}(z_j)$ s are used instead of $q^{(t)}(z_j)$ s for the distribution on the right-hand side of (31). It depends on the hierarchical structure of the model. Similarly, some $q^{(0)}(z_i)$ s may not be necessary.

B. Taylor Approximations

Although variational Bayes is a widely used general framework, its application is difficult in practice because it requires a conjugate prior. The prior distributions we have introduced are not conjugate priors. However, we have found that simple Taylor approximations make them conjugate and enable the analytical exact expectations in (31).

Here, to simplify the notation, we define the mean values of the latent variables $\boldsymbol{\eta}$, the hyper parameters $\lambda, \rho, \kappa, \beta$, and the registration parameters $\boldsymbol{\phi}_l$ over the trial distributions in the step number t of the updates of variational Bayes as $\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, \mu_{\lambda}^{(t)}, \mu_{\rho}^{(t)}, \mu_{\kappa}^{(t)}, \mu_{\beta}^{(t)}, \boldsymbol{\mu}_{\boldsymbol{\phi}_l}^{(t)}$.

Specifically, we use first-order Taylor approximations for three non-linear terms. $\mathbf{W}(\boldsymbol{\phi}_l)$ is approximated around $\boldsymbol{\phi}_l = \boldsymbol{\mu}_{\boldsymbol{\phi}_l}^{(t)}$,

$$\mathbf{W}(\boldsymbol{\phi}_l) \simeq \mathbf{W}_l^{(t)} + \sum_{k=1}^4 [\boldsymbol{\phi}_l - \boldsymbol{\mu}_{\boldsymbol{\phi}_l}^{(t)}]_k \mathbf{W}_{l,k}'^{(t)}, \quad (32)$$

where

$$\mathbf{W}_l^{(t)} \equiv \mathbf{W}(\boldsymbol{\mu}_{\boldsymbol{\phi}_l}^{(t)}), \quad (33)$$

$$\mathbf{W}_{l,k}'^{(t)} \equiv \left. \frac{\partial \mathbf{W}(\boldsymbol{\phi}_l)}{\partial \phi_{l,k}} \right|_{\boldsymbol{\phi}_l = \boldsymbol{\mu}_{\boldsymbol{\phi}_l}^{(t)}}. \quad (34)$$

Similarly, $\ln |\mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)|$ is approximated around $[\boldsymbol{\eta}, \ln \rho, \ln \kappa] = [\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, \ln \mu_{\rho}^{(t)}, \ln \mu_{\kappa}^{(t)}]$,

$$\begin{aligned} \ln |\mathbf{A}(\boldsymbol{\eta}, \rho, \kappa)| &\simeq \ln |\mathbf{A}(\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, \mu_{\rho}^{(t)}, \mu_{\kappa}^{(t)})| \\ &+ \operatorname{tr} \mathbf{A}(\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, \mu_{\rho}^{(t)}, \mu_{\kappa}^{(t)})^{-1} \left[\mu_{\rho}^{(t)} \mathbf{A}(\boldsymbol{\eta} - \boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, 1, 0) \right. \\ &\left. + (\ln \rho - \ln \mu_{\rho}^{(t)}) \mu_{\rho}^{(t)} \mathbf{A}(\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t)}, 1, 0) + (\ln \kappa - \ln \mu_{\kappa}^{(t)}) \mu_{\kappa}^{(t)} \mathbf{I} \right]. \end{aligned} \quad (35)$$

We also use a similar approximation around $[\boldsymbol{\eta}, \ln \rho, \ln \kappa] = [\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(t+1)}, \ln \mu_{\rho}^{(t)}, \ln \mu_{\kappa}^{(t)}]$. In addition, $\ln \text{logistic}(\lambda)$ is approximated around $\ln \lambda = \ln \mu_{\lambda}^{(t)}$,

$$\begin{aligned} \ln \text{logistic}(\lambda) &\simeq \ln \text{logistic}(\mu_{\lambda}^{(t)}) \\ &+ (\ln \lambda - \ln \mu_{\lambda}^{(t)}) \mu_{\lambda}^{(t)} \text{logistic}(-\mu_{\lambda}^{(t)}). \end{aligned} \quad (36)$$

C. Update Equations

The trial distributions are obtained from (30)-(32), (35), and (36), as follows:

$$q^{(\theta)}(\boldsymbol{\eta}) = \prod_{i \sim j} \text{Bernoulli}(\eta_{i,j}; \mu_{\eta_{i,j}}^{(\theta)}), \quad (37)$$

$$q^{(\theta)}(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_{\mathbf{x}}^{(\theta)}, \boldsymbol{\Sigma}_{\mathbf{x}}^{(\theta)}), \quad (38)$$

$$q^{(\theta)}(\lambda, \rho, \kappa, \beta) = \text{Gamma}(\lambda; a_{\lambda}^{(\theta)}, b_{\lambda}^{(\theta)}) \text{Gamma}(\rho; a_{\rho}^{(\theta)}, b_{\rho}^{(\theta)}) \\ \times \text{Gamma}(\kappa; a_{\kappa}^{(\theta)}, b_{\kappa}^{(\theta)}) \text{Gamma}(\beta; a_{\beta}^{(\theta)}, b_{\beta}^{(\theta)}), \quad (39)$$

$$q^{(\theta)}(\Phi) = \prod_{l=1}^L \mathcal{N}(\phi_l; \boldsymbol{\mu}_{\phi_l}^{(\theta)}, \boldsymbol{\Sigma}_{\phi_l}^{(\theta)}). \quad (40)$$

For (30) and (31), we update those distributions as follows. First, we compute $q^{(t+1)}(\boldsymbol{\eta})$ using $q^{(\theta)}(\mathbf{x}, \lambda, \rho, \kappa, \beta, \Phi)$. Second, we compute $q^{(t+1)}(\mathbf{x})$ using $q^{(t+1)}(\boldsymbol{\eta})q^{(\theta)}(\lambda, \rho, \kappa, \beta, \Phi)$. Finally, we compute $q^{(t+1)}(\lambda, \rho, \kappa, \beta)$ using $q^{(t+1)}(\mathbf{x}, \boldsymbol{\eta})q^{(\theta)}(\Phi)$ and $q^{(t+1)}(\Phi)$ using $q^{(t+1)}(\mathbf{x}, \boldsymbol{\eta})q^{(\theta)}(\lambda, \rho, \kappa, \beta)$. Here, we simply compute only the parameters of those distributions because we can compute the expectations in (31) analytically by using Taylor approximations in (32), (35), and (36). Specific update equations are described in the Appendix.

For the initial parameters of the trial distributions of $\boldsymbol{\eta}$ and \mathbf{x} , we use non-informative values,

$$\boldsymbol{\mu}_{\boldsymbol{\eta}}^{(0)} \equiv \mathbf{0}, \quad \boldsymbol{\mu}_{\mathbf{x}}^{(0)} \equiv \mathbf{0}, \quad \boldsymbol{\Sigma}_{\mathbf{x}}^{(0)} \equiv \mathbf{0}. \quad (41)$$

For the initial parameters for $\lambda, \rho, \kappa, \beta$ and Φ , we use the same values as their prior's values.

We obtain the well-approximated PM of \mathbf{x} as $\boldsymbol{\mu}_{\mathbf{x}}^{(\infty)}$. Realistically, instead of $\boldsymbol{\mu}_{\mathbf{x}}^{(\infty)}$, we use $\boldsymbol{\mu}_{\mathbf{x}}^{(t+1)}$ when the following convergence conditions hold for $\boldsymbol{\mu}_{\mathbf{x}}^{(t+1)}$ and each $\mu_{\phi_{l,k}}^{(t+1)}$,

$$\frac{1}{N_{\mathbf{x}}} \|\boldsymbol{\mu}_{\mathbf{x}}^{(t+1)} - \boldsymbol{\mu}_{\mathbf{x}}^{(t)}\|_2^2 < 10^{-4}, \\ \frac{1}{L} \sum_{l=1}^L \frac{(\mu_{\phi_{l,k}}^{(t+1)} - \mu_{\phi_{l,k}}^{(t)})^2}{[\sigma_{\phi_{l,k}}^2]} < 10^{-4} \quad (k = 1, 2, 3, 4), \quad (42)$$

where we defined $\sigma_{\phi}^2 \equiv [10^{-3}, 10^0, 10^0, 10^{-3}]$ as the scaling constant.

V. EXPERIMENTAL RESULTS

The proposed method was evaluated using five gray-scale images with a size of 40×40 pixels, as shown in Fig. 2. From each image, $L = 10$ images with a size of 10×10 pixels were created by using (1), (2) with the settings of the parameters α, Φ , and β as the following. The resolution enhancement factor α was 4. The transformation parameter Φ was randomly created according to the prior distribution in (19). The noise level parameter β was set for signal-to-noise ratios (SNR) of 20, 25, and 30 dB for each image. Samples of the created images are shown in Fig. 3.

Figure 4 shows the images estimated under SNR= 30dB. The resolution of each image appeared to be better than the corresponding observed image in Fig. 3.

Table I lists the quantitative results compared to those from the methods of bilinear interpolation, Kanemura *et al.* [6], and

Babacan *et al.* [10]. Note that we added a slight modification to these methods because they employ slightly different models. For example, the original method [10] assumes the blurring parameter γ is known, so we set γ as the mean value of the true distribution for this method. Also, we introduced a strong prior for λ in the Kanemura method [6] in contrast to the original method, because this parameter sometimes becomes negative. We evaluated the results with regard to the expectation and the standard deviation of the improvement in signal-to-noise ratio (ISNR) over 10 experiments on each image and for each SNR. ISNR is the relative PSNR defined as

$$\text{ISNR} \equiv \text{PSNR}(\hat{\mathbf{x}}; \mathbf{x}) - \text{PSNR}(\tilde{\mathbf{x}}; \mathbf{x}), \quad (43)$$

where \mathbf{x} is the true HR image, $\hat{\mathbf{x}}$ is the image estimated by the proposed method, and $\tilde{\mathbf{x}}$ is the image estimated by the compared method. A higher ISNR value means better improvement of the estimate against the estimate of the compared method. We see that the ISNRs of the proposed method were mostly higher than those of the other methods, except for the comparison with the Babacan's method in Pepper image.

Table II lists the root mean square errors (RMSE) of our method and the other methods. To evaluate the estimated registration parameters, we took the RMSEs over 50 experiments (10 experiments \times 5 images) for each noise level. Of course, a lower RMSE value means a better estimate. We see that the RMSEs of the proposed method were mostly higher than those of the other methods.

The calculation times of the proposed method was about 10 minutes on an Intel Core i7 2600 processor. The proposed method was a little slower than the method of Babacan *et al.* [10] and a little faster than the method of Kanemura *et al.* [6].

VI. DISCUSSION

With regard to the observation model, we used a linear transformation and AWGN. The use of the linear transformation model is advantageous since an arbitrary transformation matrix $\mathbf{W}(\phi_l)$ can be employed because of the Taylor approximation. The transformation matrix can be constructed by multiplying three matrices: the warping, blurring, and downsampling matrices [10]. A disadvantage of this is that sub-pixel errors might accumulate. We prefer matrix construction via a continuous function [3]. We improved the construction by introducing an elliptic theta function for the normalizing constant in (6). This normalizing constant provides fair pixel weights for both marginal and central areas of the HR image and faithfully represents the Gaussian PSF.

With regard to the HR image prior, we used a causal type of prior, which was first introduced by Kanemura *et al.* [6], [8]. The microstate energy function, or equivalently, the Hamiltonian, -based compound MRF prior of (14), offers the advantage of easy construction, but it usually has an exponential calculation cost, $\mathcal{O}(2^{N_{\eta}})$, for the normalizing constant or, equivalently, the partition function, and this is an obstacle to direct calculation of the PM solution. The MAP solution has been used in work elsewhere because it is not affected by the normalizing constant. In contrast, the introduced causal type of prior of (8) has only a polynomial

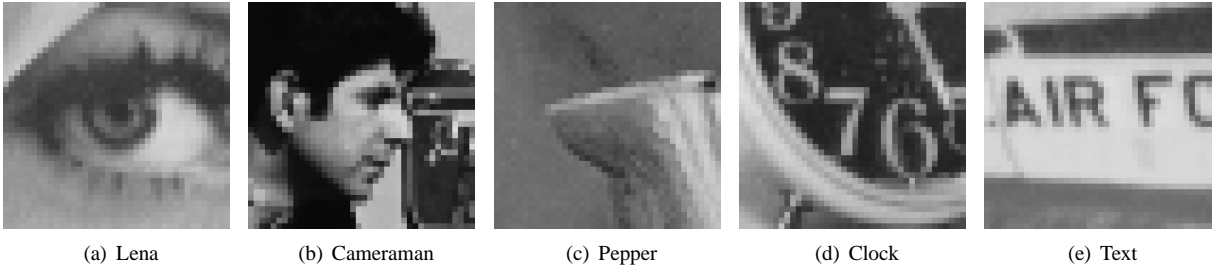


Fig. 2. Five original images used in the experiments

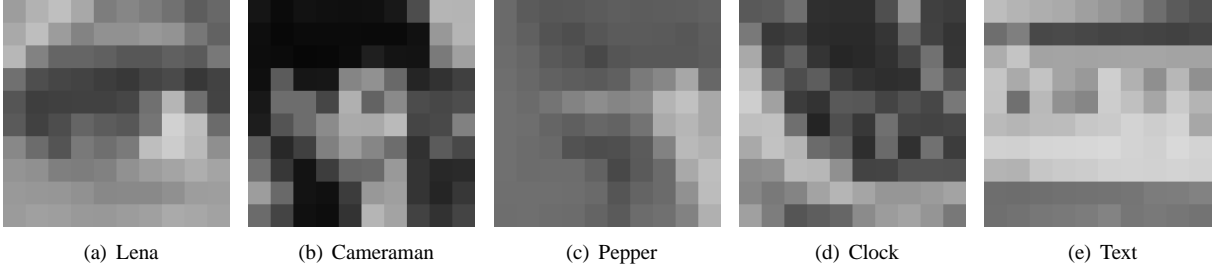


Fig. 3. Observed images when warped, blurred, downsampled by an enhancement factor of 4, and noised with SNR= 30dB AWGN

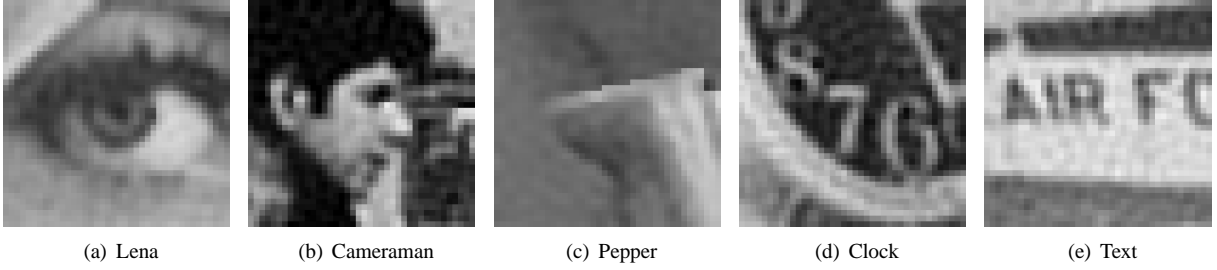


Fig. 4. Images estimated from Fig. 3 observed images

calculation cost $\mathcal{O}(N_x^3)$, which enables us to successfully apply the variational Bayes method to this problem.

With regard to the hyperparameter priors, we also improved the existing method. As the edge penalty parameter λ , Kanemura *et al.* [6] implicitly assumed $\lambda \in \mathbb{R}$, which leads to a negative λ and consequently results in an edge-strewn image. We assumed $\lambda > 0$ by setting its prior according to a gamma distribution, resulting in an appropriate inference. As the smoothness parameter ρ , they practically fixed the value of ρ with a strongly informative prior. We chose a non-informative prior for ρ . We show the box and whisker plot of the PM for each hyperparameter over 10 experiments on each image under SNR= 30dB noise in Fig. 5. As can be seen, the inferred value of the PM of ρ showed wide variation, with an approximately 10-fold maximum-to-minimum ratio, depending on the original image. This result can be interpreted as meaning it is worth inferring ρ in each HR image. Furthermore, λ and κ respectively showed approximately 2-fold and 4-fold ranges of variation. Regarding the contrast parameter κ , they assumed $\kappa \equiv 0$, which leads to $|\mathbf{A}| = 0$, and this results in an improper normalizing constant. While we assume $\kappa > 0$, which leads to a proper normalizing constant, we can consequently take the term of $\ln |\mathbf{A}|$ into account in the update equations of the variational Bayes.

With regard to the prior distribution for the blurring param-

eter γ , we used a Gaussian distribution even though γ is a positive real number. This is because we selected a simpler expression. We tried using the prior of the gamma distribution as γ , but the improvement was small. One disadvantage of this model is that a non-informative setting for this prior may lead to a nonsense result where the inferred γ is negative. Moreover, we employed a somewhat informative prior for γ . This is because the blurring parameter γ and smoothness hyperparameter ρ are somewhat complementary. This means that simultaneous estimation of γ and ρ is difficult. Tipping *et al.* [3] and Kanemura *et al.* [6] fixed ρ , and Babacan *et al.* [9] fixed γ .

With regard to the estimator, we logically derived the optimal estimator PM from the objective function of the L2-norm-based PSNR. The widely used joint MAP estimator can be considered the optimal estimator for the all-or-none type objective function,

$$\operatorname{argmax}_{\hat{z}} \langle \delta(\hat{z} - z) \rangle_{p(z|\mathbf{Y})} = \operatorname{argmax}_z p(z|\mathbf{Y}), \quad (44)$$

where δ is the Dirac delta or Kronecker delta function. Generally, this type of objective function is nonsensical for continuous variables because it is measure zero. If all the random variables in the posterior distribution are discrete, or if we can assume some smoothness of the posterior distribution,

TABLE I

PSNR OF THE PROPOSED METHOD (A HIGHER VALUE IS BETTER) AND ISNRs AGAINST THREE PREVIOUS METHODS (A HIGHER VALUE IS BETTER) FOR DIFFERENT IMAGES AND SNR LEVELS

Image	SNR [dB]	PSNR (proposed)	ISNR		
			(vs bilinear)	(vs Kanemura)	(vs Babacan)
Lena	20	29.22 ± 0.33	+5.35 ± 0.33	+0.73 ± 0.36	−0.11 ± 0.08
	25	30.69 ± 0.25	+6.78 ± 0.30	+1.04 ± 0.32	+0.11 ± 0.11
	30	32.13 ± 0.36	+8.18 ± 0.39	+1.60 ± 0.37	+0.48 ± 0.24
Cameraman	20	21.74 ± 0.19	+4.11 ± 0.20	+1.06 ± 0.36	+0.05 ± 0.07
	25	22.68 ± 0.21	+5.00 ± 0.23	+1.16 ± 0.34	−0.06 ± 0.08
	30	23.72 ± 0.38	+6.04 ± 0.39	+1.81 ± 0.17	+0.03 ± 0.06
Pepper	20	29.71 ± 0.37	+3.69 ± 0.35	+0.11 ± 0.17	+0.28 ± 0.12
	25	30.69 ± 0.28	+4.57 ± 0.27	+0.28 ± 0.32	+0.06 ± 0.16
	30	31.23 ± 0.42	+5.10 ± 0.43	+0.81 ± 0.75	−0.34 ± 0.48
Clock	20	23.27 ± 0.18	+5.36 ± 0.19	+1.49 ± 0.22	+0.11 ± 0.10
	25	24.29 ± 0.22	+6.35 ± 0.22	+1.77 ± 0.26	+0.09 ± 0.08
	30	25.49 ± 0.50	+7.53 ± 0.51	+2.49 ± 0.19	+0.32 ± 0.13
Text	20	24.67 ± 0.26	+5.83 ± 0.28	+1.57 ± 0.20	−0.10 ± 0.04
	25	25.87 ± 0.30	+7.00 ± 0.33	+1.98 ± 0.33	−0.03 ± 0.14
	30	27.26 ± 0.65	+8.37 ± 0.67	+3.03 ± 0.42	+0.18 ± 0.06

a joint MAP solution will have meaning. Instead of the L2-norm-based objective function of PSNR, the L1-norm (mean absolute error) -based PSNR is sometimes employed. In such cases, the median of the posterior distribution is generally the optimal estimator. In the case of the marginalized ML, or equivalently, type-II ML or empirical Bayes, for example, the registration parameters and other hyperparameters are firstly inferred as:

$$[\hat{\lambda}, \hat{\rho}, \hat{\kappa}, \hat{\beta}, \hat{\Phi}] \equiv \operatorname{argmax}_{\lambda, \rho, \kappa, \beta, \Phi} p(\mathbf{Y} | \lambda, \rho, \kappa, \beta, \Phi). \quad (45)$$

If these parameters have priors, such a method is called marginalized MAP. The HR image and sometimes the edge information are then inferred to as MAP,

$$\hat{x} \equiv \operatorname{argmax}_x \max_{\eta} p(x, \eta | \mathbf{Y}, \hat{\lambda}, \hat{\rho}, \hat{\kappa}, \hat{\beta}, \hat{\Phi}), \quad (46)$$

or PM. For such a two-step inference, it is difficult to calculate back the objective function.

With regard to the Taylor approximation for the transformation matrix $\mathbf{W}(\phi_l)$, we used the first-order approximation in (32) because it is more stable than the second-order approximation. This first-order approximation was proposed by Villena *et al.* [9]. The second-order approximation was proposed by Pickup *et al.* [5], and they obtained good results. We also tried the second-order approximation, but it sometimes made the algorithm unstable because it sometimes failed to produce a positive definite matrix for the covariance matrix Σ_x .

With regard to the Taylor approximation for $\ln |\mathbf{A}(\eta, \rho, \kappa)|$ and $\ln \operatorname{logistic}(\lambda)$, we introduced the first-order approximation around $[\eta, \ln \rho, \ln \kappa] = [\mu_\eta^{(t)}, \ln \mu_\rho^{(t)}, \ln \mu_\kappa^{(t)}]$ and $\ln \lambda = \ln \mu_\lambda^{(t)}$, respectively, in (35) and (36). Note that the Taylor expansion not with respect to ρ, κ, λ , but with respect to $\ln \rho, \ln \kappa, \ln \lambda$ is our key idea to solve the conjugate prior problem. Indeed, we could successfully derive the terms originating from $\ln |\mathbf{A}|$ in update equations ((52), (60), and (62) in Appendix). Kanemura *et al.* [6], [8] ignored the term of $\ln |\mathbf{A}|$ because of the high calculation cost, and this would result in less accurate inference. As for η , we implicitly assumed that η is not a binary vector but a continuous vector and did the differentiation. This assumption is based on (12). If we make another assumption

– i.e., replacement of $\eta_{i,j}$ with $\eta_{i,j}^2$ in (12) – (12) has the same meaning, but the result of the Taylor approximation will differ from the current form.

With regard to the experimental results, the proposed method outperforms the other methods in terms of the ISNR for most images and noise levels. Moreover, its estimation of the registration parameters was more accurate than the other methods were for most conditions. Therefore, we conclude the proposed method is on the whole superior to the other methods. Compared with bilinear interpolation and Kanemura's method, the superiority of the proposed method was clear. Compared with the Babacan's method, the superiority of the proposed method was rather slight. Especially, in the case of the Pepper image in 30 dB noise, the proposed method was worse than the Babacan's method. This inferiority is considered to be caused by unstable estimation of γ and ρ , where Babacan's method fixed the value of γ to the true expected value in our implementation. Intuitively, the Pepper image is smoother than the other images and has fewer edges. Therefore, this feature is considered to be less preferable for complementary parameters of γ and ρ .

With regard to the calculation cost, the proposed algorithm requires $\mathcal{O}(N_x^3)$. This calculation cost order is given by two matrix inversions: $\Sigma_x^{(t+1)}$ in (55) and \mathbf{A} in (52) and (62) (see Appendix). We found that a simple approximation such as considering all the off-diagonal elements to be zero reduces the calculation time but obviously degrades accuracy. We hope to solve this problem in our future work.

VII. CONCLUSION

In this paper, we proposed a Bayesian image super-resolution (SR) method with a causal Gaussian Markov random field (MRF) prior. We improved existing models with respect to three points: 1) the combined transformation model through a preferable normalization term using the elliptic theta function, 2) the causal Gaussian MRF model through introduction of a contrast parameter κ , which provides an effective normalizing constant including $\ln |\mathbf{A}|$, and 3) the hyperparameter prior model through application of a gamma

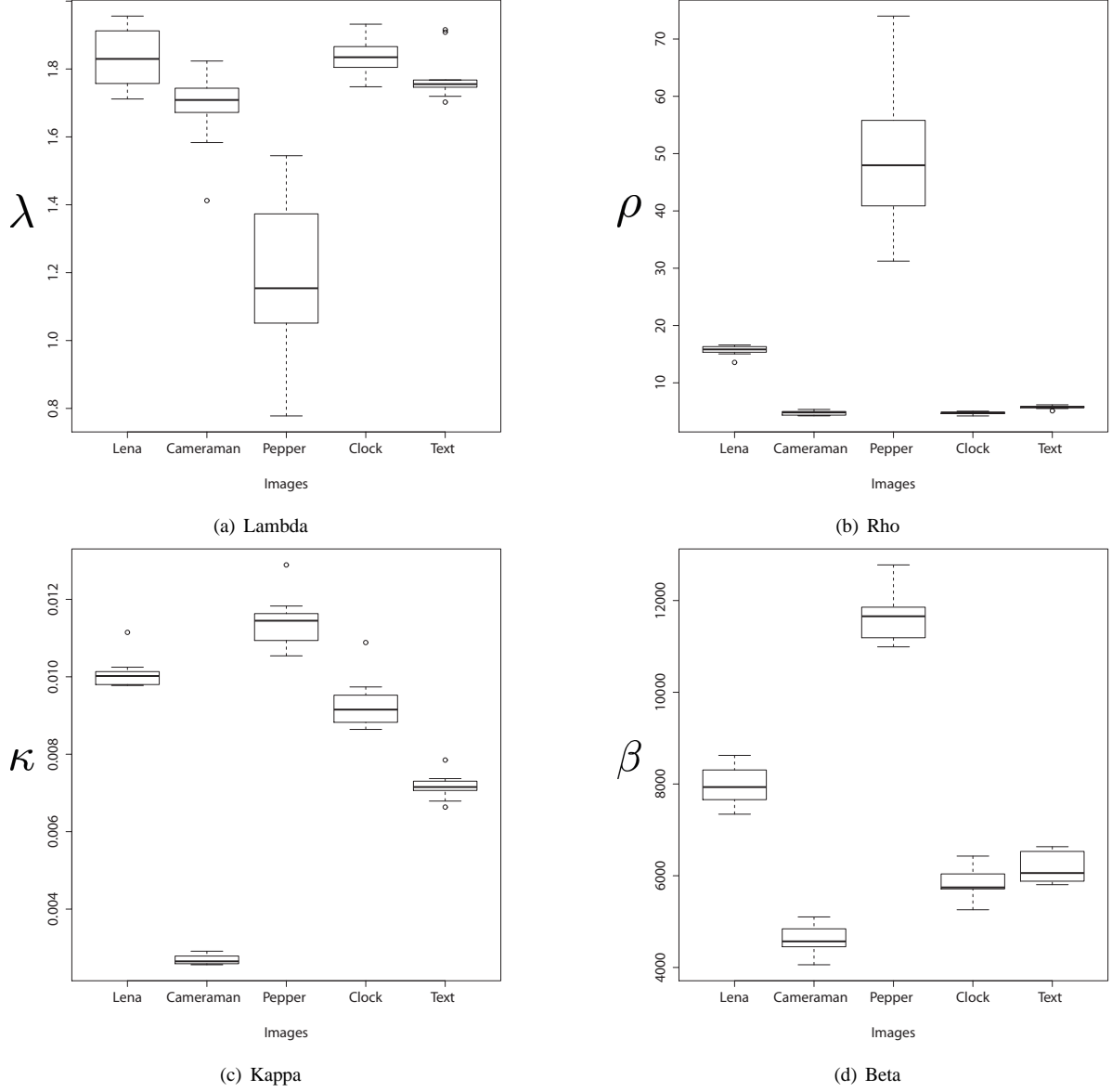


Fig. 5. Box and whisker plot of the PM for each hyperparameter, λ , ρ , κ , and β , and image under SNR= 30dB noise

TABLE II
RMSES OF REGISTRATION PARAMETERS (A LOWER VALUE IS BETTER)
FOR DIFFERENT SNR LEVELS

parameter	SNR [dB]	RMSE		
		(proposed)	(Kanemura)	(Babacan)
θ	20	0.006	0.006	0.006
	25	0.004	0.004	0.004
	30	0.002	0.003	0.003
$[\hat{o}]_h$	20	0.094	0.095	0.094
	25	0.054	0.059	0.056
	30	0.041	0.060	0.046
$[\hat{o}]_v$	20	0.074	0.073	0.076
	25	0.044	0.052	0.047
	30	0.037	0.044	0.036
γ	20	0.031	0.033	—
	25	0.025	0.030	—
	30	0.028	0.028	—

an unfavorable edge-strewn image. We then logically derived the optimal estimator, that is, not the joint maximum a posteriori (MAP) or marginalized maximum likelihood (ML) but the posterior mean (PM), from the objective function of the L2-norm (mean square error) -based peak signal-to-noise ratio (PSNR). The estimator is numerically determined by using variational Bayes. We solved the conjugate prior problem in variational Bayes by introducing three Taylor approximations. Other than these approximations, we did not use any approximations such as ignoring the term $\ln |\mathbf{A}|$. Experimental results showed that the proposed method is mostly superior to existing methods in accuracy.

APPENDIX

Here, we show the details of the variational Bayes' update equations in Section IV-C.

distribution for the edge penalty parameter λ , which prevents

The mean values of the hyperparameters $\lambda, \rho, \kappa, \beta$ over the trial distributions $q^{(t)}(\lambda, \rho, \kappa, \beta)$ are given by

$$\mu_\lambda^{(t)} = \frac{a_\lambda^{(t)}}{b_\lambda^{(t)}}, \mu_\rho^{(t)} = \frac{a_\rho^{(t)}}{b_\rho^{(t)}}, \mu_\kappa^{(t)} = \frac{a_\kappa^{(t)}}{b_\kappa^{(t)}}, \mu_\beta^{(t)} = \frac{a_\beta^{(t)}}{b_\beta^{(t)}}. \quad (47)$$

The update equation of η is given as

$$\begin{aligned} q^{(t+1)}(\eta) &\propto \exp \langle \ln p(\mathbf{z}|\mathbf{Y}) \rangle_{q^{(t)}(\mathbf{x}, \lambda, \rho, \kappa, \beta, \Phi)} \\ &\propto \exp \left(\sum_{i \sim j} \left\{ c_\lambda^{(t)} - \frac{\mu_\rho^{(t)}}{2} \text{tr} \mathbf{C}_x^{(t)} \mathbf{M}_{i,j} \right\} \eta_{i,j} \right. \\ &\quad \left. + \frac{1}{2} \langle \ln |\mathbf{A}(\eta, \rho, \kappa)| \rangle_{q^{(t)}(\rho, \kappa)} \right), \end{aligned} \quad (48)$$

where

$$\mathbf{C}_x^{(t)} \equiv \mu_x^{(t)} [\mu_x^{(t)}]^\top + \Sigma_x^{(t)}, \quad (49)$$

$$[\mathbf{M}_{i,j}]_{k,l} \equiv \begin{cases} +1, & (k,l) = (i,i) \text{ or } (j,j), \\ -1, & (k,l) = (i,j) \text{ or } (j,i), \\ 0, & \text{otherwise.} \end{cases} \quad (50)$$

Using the Taylor approximation of (35), we obtain the distribution of (37) at step $t+1$ with the parameter

$$\mu_{\eta_{i,j}}^{(t+1)} = \text{logistic} \left(\mu_\lambda^{(t)} + \frac{1}{2} \mu_\rho^{(t)} C_{\eta_{i,j}}^{(t)} \right), \quad (51)$$

where

$$C_{\eta_{i,j}}^{(t)} \equiv \text{tr} \left[\left(\mathbf{A}(\mu_\eta^{(t)}, \mu_\rho^{(t)}, \mu_\kappa^{(t)})^{-1} - \mathbf{C}_x^{(t)} \right) \mathbf{M}_{i,j} \right]. \quad (52)$$

The update equation of \mathbf{x} is given as

$$\begin{aligned} q^{(t+1)}(\mathbf{x}) &\propto \exp \langle \ln p(\mathbf{z}|\mathbf{Y}) \rangle_{q^{(t+1)}(\eta) q^{(t)}(\lambda, \rho, \kappa, \beta, \Phi)} \\ &\propto \exp \left(-\frac{1}{2} \left\{ \mathbf{x}^\top \mathbf{A}(\mu_\eta^{(t+1)}, \mu_\rho^{(t)}, \mu_\kappa^{(t)}) \mathbf{x} \right. \right. \\ &\quad \left. \left. + \mu_\beta^{(t)} \sum_{l=1}^L \langle \|\mathbf{W}(\phi_l) \mathbf{x} - \mathbf{y}_l\|_2^2 \rangle_{q^{(t)}(\phi_l)} \right\} \right). \end{aligned} \quad (53)$$

It becomes a Gaussian distribution. Using the Taylor approximation (32), we obtain the distribution of (38) at step $t+1$ with the parameters

$$\mu_x^{(t+1)} = \Sigma_x^{(t+1)} \left[\mu_\beta^{(t)} \sum_{l=1}^L \mathbf{y}_l^\top \mathbf{W}_l^{(t)} \right]^\top, \quad (54)$$

$$\Sigma_x^{(t+1)} = \left[\mathbf{A}(\mu_\eta^{(t+1)}, \mu_\rho^{(t)}, \mu_\kappa^{(t)}) + \mu_\beta^{(t)} \sum_{l=1}^L \mathbf{C}_{\mathbf{W}_l}^{(t)} \right]^{-1}, \quad (55)$$

where

$$\mathbf{C}_{\mathbf{W}_l}^{(t)} \equiv [\mathbf{W}_l^{(t)}]^\top \mathbf{W}_l^{(t)} + \sum_{k,k'} [\Sigma_{\phi_l}^{(t)}]_{k,k'} [\mathbf{W}_{l,k}^{(t)}]^\top \mathbf{W}_{l,k'}^{(t)}. \quad (56)$$

The update equation of $\lambda, \rho, \kappa, \beta$ is given as

$$\begin{aligned} q^{(t+1)}(\lambda, \rho, \kappa, \beta) &\propto \exp \langle \ln p(\mathbf{z}|\mathbf{Y}) \rangle_{q^{(t+1)}(\mathbf{x}, \eta) q^{(t)}(\Phi)} \\ &\propto \lambda^{a_\lambda^{(0)}-1} \rho^{a_\rho^{(0)}-1} \kappa^{a_\kappa^{(0)}-1} \beta^{a_\beta^{(0)}+\frac{1}{2}LN_y} \exp \left(\right. \\ &\quad - \left\{ b_\lambda^{(0)} + \sum_{i \sim j} (1 - \mu_{\eta_{i,j}}^{(t+1)}) \right\} \lambda \\ &\quad - \left\{ b_\rho^{(0)} + \frac{1}{2} \text{tr} \mathbf{C}_x^{(t+1)} \mathbf{A}(\mu_\eta^{(t+1)}, 1, 0) \right\} \rho \\ &\quad - \left\{ b_\kappa^{(0)} + \frac{1}{2} \text{tr} \mathbf{C}_x^{(t+1)} \right\} \kappa \\ &\quad - \left\{ b_\beta^{(0)} + \frac{1}{2} \sum_{l=1}^L \left\langle \text{tr} \mathbf{C}_x^{(t+1)} \mathbf{W}(\phi_l)^\top \mathbf{W}(\phi_l) \right. \right. \\ &\quad \left. \left. - 2 \mathbf{y}_l^\top \mathbf{W}(\phi_l) \mu_x^{(t+1)} + \mathbf{y}_l^\top \mathbf{y}_l \right\rangle_{q^{(t)}(\phi_l)} \right\} \beta \\ &\quad \left. + \frac{1}{2} \langle \ln |\mathbf{A}(\eta, \rho, \kappa)| \rangle_{q^{(t+1)}(\eta)} + N_\eta \ln \text{logistic}(\lambda) \right). \end{aligned} \quad (57)$$

Using Taylor approximations (32), (35), and (36), we obtain the distribution of (39) at step $t+1$ with parameters

$$a_\lambda^{(t+1)} = a_\lambda^{(0)} + N_\eta \mu_\lambda^{(t)} \text{logistic}(-\mu_\lambda^{(t)}), \quad (58)$$

$$b_\lambda^{(t+1)} = b_\lambda^{(0)} + \sum_{i \sim j} (1 - \mu_{\eta_{i,j}}^{(t+1)}), \quad (59)$$

$$a_\rho^{(t+1)} = a_\rho^{(0)} + \frac{\mu_\rho^{(t)}}{2} \text{tr} \mathbf{A}(\mu_\eta^{(t+1)}, \mu_\rho^{(t)}, \mu_\kappa^{(t)})^{-1} \mathbf{A}(\mu_\eta^{(t+1)}, 1, 0) \quad (60)$$

$$b_\rho^{(t+1)} = b_\rho^{(0)} + \frac{1}{2} \text{tr} \mathbf{C}_x^{(t+1)} \mathbf{A}(\mu_\eta^{(t+1)}, 1, 0), \quad (61)$$

$$a_\kappa^{(t+1)} = a_\kappa^{(0)} + \frac{\mu_\kappa^{(t)}}{2} \text{tr} \mathbf{A}(\mu_\eta^{(t+1)}, \mu_\rho^{(t)}, \mu_\kappa^{(t)})^{-1} \quad (62)$$

$$b_\kappa^{(t+1)} = b_\kappa^{(0)} + \frac{1}{2} \text{tr} \mathbf{C}_x^{(t+1)}, \quad (63)$$

$$a_\beta^{(t+1)} = a_\beta^{(0)} + \frac{1}{2} LN_y, \quad (64)$$

$$b_\beta^{(t+1)} = b_\beta^{(0)} + \frac{1}{2} \sum_{l=1}^L \left(\text{tr} \mathbf{C}_x^{(t+1)} \mathbf{C}_{\mathbf{W}_l}^{(t)} - 2 \mathbf{y}_l^\top \mathbf{W}_l^{(t)} \mu_x^{(t+1)} + \mathbf{y}_l^\top \mathbf{y}_l \right). \quad (65)$$

The update equation of Φ is given as

$$\begin{aligned} q^{(t+1)}(\Phi) &\propto \exp \langle \ln p(\mathbf{z}|\mathbf{Y}) \rangle_{q^{(t+1)}(\mathbf{x}, \eta) q^{(t)}(\lambda, \rho, \kappa, \beta)} \\ &\propto \exp \left(-\frac{1}{2} \sum_{l=1}^L \left\{ [\phi_l - \mu_{\phi_l}^{(0)}]^\top [\Sigma_{\phi_l}^{(0)}]^{-1} [\phi_l - \mu_{\phi_l}^{(0)}] \right. \right. \\ &\quad \left. \left. + \mu_\beta^{(t)} \left\{ \text{tr} \mathbf{C}_x^{(t+1)} \mathbf{W}(\phi_l)^\top \mathbf{W}(\phi_l) - 2 \mathbf{y}_l^\top \mathbf{W}(\phi_l) \mu_x^{(t+1)} \right\} \right\} \right). \end{aligned} \quad (66)$$

Using the Taylor approximation (32), we obtain the distribu-

tion of (40) at step $t + 1$ with parameters

$$\mu_{\phi_l}^{(t+1)} = \Sigma_{\phi_l}^{(t+1)} \left[[\Sigma_{\phi_l}^{(0)}]^{-1} \mu_{\phi_l}^{(0)} + \mu_{\beta}^{(t)} [C_{\phi_l}^{\prime\prime(t+1)} \mu_{\phi_l}^{(t)} - C_{\phi_l}^{\prime(t+1)}] \right], \quad (67)$$

$$\Sigma_{\phi_l}^{(t+1)} = \left[[\Sigma_{\phi_l}^{(0)}]^{-1} + \mu_{\beta}^{(t)} C_{\phi_l}^{\prime\prime(t+1)} \right]^{-1}, \quad (68)$$

where

$$[C_{\phi_l}^{\prime(t+1)}]_k \equiv \frac{1}{2} \text{tr} C_x^{(t+1)} \left[[W_l^{(t)}]^\top W_{l,k}^{\prime(t)} + [W_{l,k}^{\prime(t)}]^\top W_l^{(t)} \right] - y_l^\top W_{l,k}^{\prime(t)} \mu_x^{(t+1)}, \quad (69)$$

$$[C_{\phi_l}^{\prime\prime(t+1)}]_{k,k'} \equiv \text{tr} C_x^{(t+1)} [W_{l,k}^{\prime(t)}]^\top W_{l,k'}^{\prime(t)}. \quad (70)$$

REFERENCES

- [1] R. Tsai and T. Huang, "Multiframe image restoration and registration," in *Advances in computer vision and image processing*, vol. 1, no. 2, JAI Press Inc., pp. 317–339, Greenwich, CT, 1984.
- [2] R. Hardie, K. Barnard, and E. Armstrong, "Joint Map registration and high resolution image estimation using a sequence of undersampled images," *IEEE Trans. Image Process.*, vol. 6, no. 12, pp. 1621–1633, 1997.
- [3] M. E. Tipping, and C. M. Bishop, "Bayesian image super-resolution," in *Advances in NIPS 15*, MIT Press, pp. 1279–1286, 2003.
- [4] R. Molina, J. Mateos, A. K. Katsaggelos, and M. Vega, "Bayesian multichannel image restoration using compound Gauss-Markov random fields," *IEEE Transactions on Image Processing*, vol. 12, no. 12, pp. 1642–1654, 2003.
- [5] L. C. Pickup, D. P. Capel, S. J. Roberts, and A. Zisserman, "Bayesian Image Super-Resolution, continued," in *Advances in NIPS 19*, MIT Press, 2007.
- [6] A. Kanemura, S. Maeda, and S. Ishii, "Hyperparameter Estimation in Bayesian Image Superresolution with a Compound Markov Random Field Prior," *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 181–186, 2007.
- [7] P. Vandewalle, L. Sbaiz, J. Vandewalle, and M. Vetterli, "Super-resolution from unregistered and totally aliased signals using subspace methods," *IEEE Trans. Signal Processing*, vol. 55, No. 7, Part 2, pp. 3687–3703, 2007.
- [8] A. Kanemura, S. Maeda, and S. Ishii, "Superresolution with compound Markov random fields via the variational EM algorithm," *Neural Networks*, vol. 22, pp. 1025–1034, 2009.
- [9] S. Villena, M. Vega, S.D. Babacan, R. Molina, and A. K. Katsaggelos, "Image Prior Combination in Super-resolution Image Registration & Reconstruction," *IEEE International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 355–360, 2010.
- [10] S. D. Babacan, R. Molina, A. K. Katsaggelos, "Variational Bayesian Super Resolution," *IEEE Transactions on Image Processing*, vol. 20, no. 4, pp. 984–999, 2011.
- [11] S. Borman and R. L. Stevenson, "Spatial resolution enhancement of low-resolution image sequences a comprehensive review with directions for future research," Department of Electrical Engineering, University of Notre Dame, Tech. Rep., 1998.
- [12] S. C. Park, M. K. Park, M. G. Kang, "Super-resolution image reconstruction: a technical overview," *Signal Processing Magazine, IEEE*, vol. 20, no.3, pp. 21–36, 2003.
- [13] S. Farsiu, D. Robinson, M. Elad, and P. Milanfar, "Advances and challenges in super-resolution," *Int. J. Imag. Syst. Technol.*, vol. 14, no. 2, pp. 47–57, 2004.
- [14] M. Ng, T. Chan, M. G. Kang, and P. Milanfar, "Special Issue on Super-resolution Imaging: Analysis, Algorithms, and Applications," *EURASIP Journal on Applied Signal Processing*, 2006.
- [15] A.K. Katsaggelos, R. Molina, and J. Mateos, "Super Resolution of Images and Video", Synthesis Lectures on Image, Video, and Multimedia Processing, Morgan & Claypool, 2007.
- [16] P. Milanfar, Ed., "Super-Resolution Imaging," CRC Press, 2010.
- [17] R. Molina, M. Vega, J. Abad, and A. K. Katsaggelos, "Parameter Estimation in Bayesian High-Resolution Image Reconstruction With Multisensors," *IEEE Transactions on Image Processing*, vol. 12, no. 12, pp. 1655–1667, 2003.
- [18] S. Geman, and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. PAMI-6, no. 6, pp. 721–741, 1984.
- [19] R. Chellappa, and A.K. Jain, Eds., "Markov random fields: Theory and application", Academic Press, Boston 1993.
- [20] F. -C. Jeng, and J. W. Woods, "Compound Gauss-Markov random fields for image estimation", *IEEE Transactions on Signal Processing*, vol. 39, no. 3, pp. 683–697, 1991.
- [21] H. Attias, "Inferring parameters and structure of latent variable models by variational Bayes," in *Proc. UAI*, San Francisco, CA, pp. 21–30, Morgan Kaufmann, 1999.
- [22] F. -C. Jeng, and J. W. Woods, "Simulated annealing in compound Gaussian random fields," *IEEE Transactions on Information Theory*, vol. 36, no. 1, pp. 94–107, 1990.
- [23] C. M. Bishop, D. Spiegelhalter, J. Winn "VIBES: A Variational Inference Engine for Bayesian Networks," in *Advances in NIPS 15*, MIT Press, pp. 777–784, 2003.